

### The case $N = M - 1$

In this case the image of  $\mathbf{f}$  is a surface in  $\mathbb{R}^M$ . If we consider the tangent plane of this surface at  $\mathbf{x}_0$ , then clearly

$$\left\{ \mathbf{y}_0 + t \frac{\partial \mathbf{f}}{\partial x_1}(\mathbf{x}_0) \mid t \in \mathbb{R} \right\}, \dots, \left\{ \mathbf{y}_0 + t \frac{\partial \mathbf{f}}{\partial x_N}(\mathbf{x}_0) \mid t \in \mathbb{R} \right\} \quad (1)$$

are  $N$  lines tangent to this surface at  $\mathbf{x}_0$ , which means the surface should be given through

$$\mathbf{n} \cdot (\mathbf{y} - \mathbf{y}_0) = 0 \quad (2)$$

where the vector  $\mathbf{n} \in \mathbb{R}^M$  satisfies

$$\mathbf{n} \cdot \frac{\partial \mathbf{f}}{\partial x_i}(\mathbf{x}_0) = 0 \quad (3)$$

for all  $i = 1, 2, \dots, N$ . This vector  $\mathbf{n}$  is called the “normal direction” of the surface at  $\mathbf{x}_0$ .

**Exercise 1.** Prove that  $\mathbf{n}$  is parallel to the following vector:

$$\det \begin{pmatrix} \mathbf{e}_1 & \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_N} \\ \vdots & \vdots & \vdots \\ \mathbf{e}_M & \frac{\partial f_M}{\partial x_1} & \frac{\partial f_M}{\partial x_N} \end{pmatrix} \quad (4)$$

where  $\mathbf{e}_1, \dots, \mathbf{e}_M$  are the coordinate directions in  $\mathbb{R}^M$ . This determinant is formal, for example

$$\det \begin{pmatrix} \mathbf{e}_1 & 3 \\ \mathbf{e}_2 & 2 \end{pmatrix} = 2 \mathbf{e}_1 - 3 \mathbf{e}_2. \quad (5)$$

**Example 1.** Consider the function

$$\mathbf{f}(t) := \begin{pmatrix} \cos t \\ \sin t \end{pmatrix}. \quad (6)$$

We calculate

$$\det \begin{pmatrix} \mathbf{e}_1 & f_1'(t) \\ \mathbf{e}_2 & f_2'(t) \end{pmatrix} = (\cos t) \mathbf{e}_1 + (\sin t) \mathbf{e}_2. \quad (7)$$

Now if we denote  $x = \cos t$ ,  $y = \sin t$ , then the vector  $\mathbf{n}$  at  $(x_0, y_0) = (\cos t_0, \sin t_0)$  is  $x_0 \mathbf{e}_1 + y_0 \mathbf{e}_2 = \begin{pmatrix} x_0 \\ y_0 \end{pmatrix}$ , the normal direction to the tangent.

**Exercise 2.** Consider the function

$$\mathbf{f}(\theta, \varphi) = \begin{pmatrix} \cos \theta \cos \varphi \\ \sin \theta \cos \varphi \\ \sin \varphi \end{pmatrix}. \quad (8)$$

Find the normal direction to its tangent plane at  $\mathbf{x}_0 := \mathbf{f}(\theta_0, \varphi_0)$  and give geometric explanation of your result.

### The case $M = 1$

Recall that in this case our function takes the form

$$f(x_1, \dots, x_N) \quad (9)$$

and the matrix representation of the differential  $Df(\mathbf{x}_0)$  is a vector, which we will call the “gradient” of  $f$  at  $\mathbf{x}_0$ .

$$(\text{grad } f)(\mathbf{x}_0) := \begin{pmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}_0) \\ \vdots \\ \frac{\partial f}{\partial x_N}(\mathbf{x}_0) \end{pmatrix}. \quad (10)$$

**Remark 2.** The case  $M = 1$  can be related to this case as a special case as follows: Given  $f: \mathbb{R}^N \mapsto \mathbb{R}$ , we define a new function  $\mathbf{g}: \mathbb{R}^N \mapsto \mathbb{R}^{N+1}$  as follows:

$$\mathbf{g}(x_1, \dots, x_N) = \begin{pmatrix} x_1 \\ \vdots \\ x_N \\ f(x_1, \dots, x_N) \end{pmatrix}. \quad (11)$$

However this identification is not very helpful for us in understanding the relation between the function and the gradient.

**Remark 3.** To avoid much confusion in future mathematics study, it is important to treat the linear function  $Df(\mathbf{x}_0)$  and its matrix/vector representation  $(\text{grad } f)(\mathbf{x}_0)$  as different objects.

**Remark 4.** Often the notation  $\nabla f$  is used for  $\text{grad } f$ . However getting too used to this notation ( $\nabla f$ ) will cause much difficulty in differential geometry.

**Remark 5.** We see that even when  $f$  is not differentiable, we may still be able to define the gradient vector – we only need the existence of all partial derivatives!

To understand the geometric meaning of the gradient  $\text{grad } f$  (not the differential  $Df!$ ), we consider the graph of  $f$ , as a subset of  $\mathbb{R}^{N+1}$ :

$$x_{N+1} = f(x_1, \dots, x_N). \quad (12)$$

**Lemma 6.** Consider the surface  $f(x_1, \dots, x_N) = c$  for some  $c \in \mathbb{R}$ . Then the gradient vector  $\text{grad } f$  is perpendicular to this surface.

**Proof.** Consider any curve  $(x_1(t), \dots, x_N(t))$  on the surface, that is

$$f(x_1(t), \dots, x_N(t)) = c. \quad (13)$$

Taking  $\frac{d}{dt}$  we have, thanks to the chain rule:

$$\frac{\partial f}{\partial x_1}(x_1(t), \dots, x_N(t)) x_1'(t) + \dots + \frac{\partial f}{\partial x_N}(x_1(t), \dots, x_N(t)) x_N'(t) = 0 \quad (14)$$

which gives

$$(\text{grad } f)(x_1(t), \dots, x_N(t)) \cdot \mathbf{x}'(t) = 0 \quad (15)$$

and conclusion follows.  $\square$

**Remark 7.** The justification of the claim:  $f(x_1, \dots, x_N) = c$  is a surface, needs Implicit function theorem, which we will discuss next week.

## Properties and applications

### Mean Value Theorem

**Theorem 8. (MVT)** Let  $f: \mathbb{R}^N \mapsto \mathbb{R}$ . Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$  and assume  $f$  is differentiable at every point on the line segment  $S := \{t\mathbf{x} + (1-t)\mathbf{y} \mid t \in [0, 1]\}$ . Then there is  $\boldsymbol{\xi} \in S$  such that

$$f(\mathbf{x}) - f(\mathbf{y}) = Df(\boldsymbol{\xi})(\mathbf{x} - \mathbf{y}). \quad (16)$$

**Proof.** Consider  $g(t) := f(t\mathbf{x} + (1-t)\mathbf{y})$ . □

**Exercise 3.** Fill in the details of the proof. In particular, why is  $g$  differentiable?

**Corollary 9.** Let  $f: E \mapsto \mathbb{R}$  be differentiable and  $E$  be convex. Then for any  $\mathbf{x}, \mathbf{y} \in E$ , there is  $\boldsymbol{\xi} \in E$  such that

$$f(\mathbf{x}) - f(\mathbf{y}) = Df(\boldsymbol{\xi})(\mathbf{x} - \mathbf{y}). \quad (17)$$

**Problem 1.** Critique the following claim:

Let  $\mathbf{f}: \mathbb{R}^N \mapsto \mathbb{R}^M$ . Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$  and assume  $\mathbf{f}$  is differentiable at every point on the line segment  $S := \{t\mathbf{x} + (1-t)\mathbf{y} \mid t \in [0, 1]\}$ . Then there is  $\boldsymbol{\xi} \in S$  such that

$$\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y}) = D\mathbf{f}(\boldsymbol{\xi})(\mathbf{x} - \mathbf{y}). \quad (18)$$

If it is correct, prove it; Otherwise find a counter-example. (Hint: You may want to consider the function  $(\cos t, \sin t)$ ).

**Problem 2.** Critique the following claim:

Let  $\mathbf{f}: \mathbb{R}^N \mapsto \mathbb{R}^M$ . Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$  and assume  $\mathbf{f}$  is differentiable at every point on the line segment  $S := \{t\mathbf{x} + (1-t)\mathbf{y} \mid t \in [0, 1]\}$ . Then

$$\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})\| \leq M \|\mathbf{x} - \mathbf{y}\| \quad (19)$$

where  $M$  can be taken as

$$M = \sup_{\mathbf{z} \in S} \left( \sum_{i,j} a_{i,j}^2(\mathbf{z}) \right)^{1/2} \quad (20)$$

with  $A = (a_{i,j})$  the Jacobian matrix of  $\mathbf{f}$ .

If it is correct, prove it; Otherwise find a counter-example.

**Question 10.** Hadamard's theorem?

### A first look at optimization theory

In this subsection we focus on the case  $f: E \subseteq \mathbb{R}^N \mapsto \mathbb{R}$ .

#### Local and global optima

**Definition 11.** Let  $f: E \subseteq \mathbb{R}^N \mapsto \mathbb{R}$ .  $\mathbf{x}_0 \in E$  is said to be

- a global maximum of  $f$  if

$$\forall \mathbf{x} \in E, \quad f(\mathbf{x}_0) \geq f(\mathbf{x}). \quad (21)$$

- a local maximum of  $f$  if there is  $r > 0$  such that

$$\forall \mathbf{x} \in B(\mathbf{x}_0, r) \cap E, \quad f(\mathbf{x}_0) \geq f(\mathbf{x}). \quad (22)$$

- a global minimum of  $f$  if

$$\forall \mathbf{x} \in E, \quad f(\mathbf{x}_0) \leq f(\mathbf{x}). \quad (23)$$

- a local minimum of  $f$  if there is  $r > 0$  such that

$$\forall \mathbf{x} \in B(\mathbf{x}_0, r) \cap E, \quad f(\mathbf{x}_0) \leq f(\mathbf{x}). \quad (24)$$

**Theorem 12.** If  $\mathbf{x}_0$  is a global maximum/minimum, then it is a local maximum/minimum.

**Exercise 4.** Prove the above theorem.

**Theorem 13.** Let  $f: E \subseteq \mathbb{R}^N \mapsto \mathbb{R}$  be differentiable and  $\mathbf{x}_0 \in E^\circ$  is a local maximum or minimum, then  $(\text{grad } f)(\mathbf{x}_0) = \mathbf{0}$ .

**Proof.** Wlog  $\mathbf{x}_0$  is a local maximum. Assume the contrary, that is  $\mathbf{v} := -(\text{grad } f)(\mathbf{x}_0) \neq \mathbf{0}$ . Then we have

$$\frac{\partial f}{\partial \mathbf{v}}(\mathbf{x}_0) = (Df)(\mathbf{x}_0)(\mathbf{v}) = (\text{grad } f)(\mathbf{x}_0) \cdot (-\mathbf{v}) = -\|\mathbf{v}\|^2 < 0. \quad (25)$$

Now by definition of directional derivative:

$$\frac{\partial f}{\partial \mathbf{v}}(\mathbf{x}_0) := \lim_{h \rightarrow 0} \frac{f(\mathbf{x}_0 + h\mathbf{v}) - f(\mathbf{x}_0)}{h} \quad (26)$$

we see that there is  $\delta > 0$  such that for all  $0 < h < \delta$ ,

$$f(\mathbf{x}_0 + h\mathbf{v}) < f(\mathbf{x}_0). \quad (27)$$

This contradicts  $\mathbf{x}_0$  being local maximum. □

**Example 14.** Consider  $f(x, y) = \sin x \sin y \sin(x + y)$ . Find its maximum/minimum on

$$E := \{(x, y) \mid x \geq 0, y \geq 0, x + y \leq \pi\} \quad (28)$$

**Solution.** As  $f$  is continuous and  $E$  is closed and bounded and thus compact, we know that  $f$  reaches its maximum and minimum on  $E$ . There are two cases for  $\mathbf{x}_0$  which is maximum/minimum: Either  $\mathbf{x}_0 \in E^\circ$  and thus satisfy  $(\text{grad } f)(\mathbf{x}_0) = \mathbf{0}$ , or  $\mathbf{x}_0 \in \partial E$ .

We first look for local maximum/minimum in the interior of the domain:

$$E^\circ := \{(x, y) \mid x > 0, y > 0, x + y < \pi\}. \quad (29)$$

Taking partial derivative we obtain

$$\frac{\partial f}{\partial x} = \cos x \sin y \sin(x + y) + \sin x \sin y \cos(x + y) = \sin(2x + y) \sin y \quad (30)$$

$$\frac{\partial f}{\partial y} = \sin(x + 2y) \sin x. \quad (31)$$

Setting both to 0 we have the following possible cases:

- $\sin(2x + y) = \sin(x + 2y) = 0$ : We must have  $x + 2y = y + 2x = \pi \implies x = y = \frac{\pi}{3}$ .
- $\sin y = \sin(x + 2y) = 0$ : There is no solution in the interior.
- $\sin(2x + y) = \sin x = 0$ : There is no solution in the interior.
- $\sin x = \sin y = 0$ : There is no solution in the interior.

Summarizing, we see that in the interior the only candidate is  $(\frac{\pi}{3}, \frac{\pi}{3})$  which gives  $f(\frac{\pi}{3}, \frac{\pi}{3}) = -\frac{3\sqrt{3}}{8}$ . On the other hand we have

$$\forall (x, y) \in \partial E, \quad f(x, y) = 0. \quad (32)$$

Thus the maximizer of  $f$  is any point on the boundary, with maximum 0; The minimizer of  $f$  is  $(\frac{\pi}{3}, \frac{\pi}{3})$  with minimum  $-\frac{3\sqrt{3}}{8}$ .

**Question 15. (Pareto optimal)** Consider functions  $f_1, \dots, f_M: \mathbb{R}^N \mapsto \mathbb{R}$ . A point  $\mathbf{x}_0 \in \mathbb{R}^N$  is called a “local Pareto maximizer” if there is  $r > 0$  such that for all  $\mathbf{x} \in B(\mathbf{x}_0, r)$ , if there is  $i \in \{1, 2, \dots, M\}$  such that  $f_i(\mathbf{x}) > f_i(\mathbf{x}_0)$ , then there must be another  $j \in \{1, 2, \dots, M\}$  such that  $f_j(\mathbf{x}_0) > f_j(\mathbf{x})$ . Explore the possibility of getting a necessary condition in terms of the Jacobian matrix of  $\mathbf{f} := \begin{pmatrix} f_1 \\ \vdots \\ f_M \end{pmatrix}$  for  $\mathbf{x}_0 \in \mathbb{R}^N$  to be a local Pareto maximizer.

## Convex optimization

**Definition 16. (Convex function)** A function  $f: A \subseteq \mathbb{R}^N \mapsto \mathbb{R}$  is convex if and only if

i.  $A$  is a convex set.

ii.  $\forall \mathbf{x}, \mathbf{y} \in A$ , and  $\forall t \in [0, 1]$ ,  $f(t\mathbf{x} + (1-t)\mathbf{y}) \leq t f(\mathbf{x}) + (1-t) f(\mathbf{y})$ .

**Theorem 17.** If  $f$  is convex, then any of its local minimum is also global.

**Proof.** Assume the contrary. Then there must be two local minimizers  $\mathbf{x}, \mathbf{y}$  such that  $f(\mathbf{x}) \neq f(\mathbf{y})$ . Wlog assume  $f(\mathbf{x}) < f(\mathbf{y})$ . Now for any  $r > 0$ , take  $t \in (0, 1)$  such that

$$t > 1 - \frac{r}{\|\mathbf{x} - \mathbf{y}\|}. \quad (33)$$

Then we have

$$\|(t\mathbf{x} + (1-t)\mathbf{y}) - \mathbf{x}\| = (1-t)\|\mathbf{x} - \mathbf{y}\| < r \implies t\mathbf{x} + (1-t)\mathbf{y} \in B(\mathbf{x}, r). \quad (34)$$

Now by convexity of  $f$  we have

$$f(t\mathbf{x} + (1-t)\mathbf{y}) \leq t f(\mathbf{x}) + (1-t) f(\mathbf{y}) < t f(\mathbf{x}) + (1-t) f(\mathbf{x}) = f(\mathbf{x}), \quad (35)$$

contradicting the assumption that  $\mathbf{x}$  is a local minimizer.  $\square$

**Theorem 18.** Let  $A$  be open and  $f: A \subseteq \mathbb{R}^N \mapsto \mathbb{R}$  be differentiable. Then  $f$  is convex if and only if for all  $\mathbf{x}, \mathbf{y} \in A$ ,

$$(Df(\mathbf{x}) - Df(\mathbf{y}))(\mathbf{x} - \mathbf{y}) \geq 0. \quad (36)$$

**Remark 19.** Such  $Df$  is said to be a “monotone”.

**Proof.**

- If. Take any  $\mathbf{x}, \mathbf{y} \in A$ . Define

$$g(t) := f(t\mathbf{x} + (1-t)\mathbf{y}). \quad (37)$$

All we need to show is for all  $t \in (0, 1)$ ,

$$g(t) \leq (1-t)g(0) + tg(1). \quad (38)$$

Now fix any  $t \in (0, 1)$ . By the one dimensional mean value theorem, we have

$$g(t) - g(0) = g'(\xi_1)t, \quad g(1) - g(t) = g'(\xi_2)(1-t) \quad (39)$$

where  $\xi_1 \in (0, t)$ ,  $\xi_2 \in (t, 1)$ . Then we have

$$(1-t)[g(t) - g(0)] + t[g(1) - g(t)] = t(1-t)[g'(\xi_1) - g'(\xi_2)]. \quad (40)$$

So all we need to show is  $g'(\xi_2) - g'(\xi_1) \geq 0$ . To do this, apply the chain rule to  $g$ :

$$g'(t) = \frac{df(t\mathbf{x} + (1-t)\mathbf{y})}{dt} = Df(t\mathbf{x} + (1-t)\mathbf{y})(\mathbf{x} - \mathbf{y}). \quad (41)$$

Now, denoting  $\mathbf{x}_i := \xi_i\mathbf{x} + (1-\xi_i)\mathbf{y}$

$$g'(\xi_2) - g'(\xi_1) = [Df(\mathbf{x}_2) - Df(\mathbf{x}_1)](\mathbf{x} - \mathbf{y}) = (\xi_2 - \xi_1)[Df(\mathbf{x}_2) - Df(\mathbf{x}_1)](\mathbf{x}_2 - \mathbf{x}_1) \quad (42)$$

which  $\geq 0$  due to  $\xi_2 > \xi_1$ .

- Only if. Take any  $\mathbf{x}, \mathbf{y} \in A$ . Define

$$g(t) := f(t\mathbf{x} + (1-t)\mathbf{y}). \quad (43)$$

Then  $g$  is convex (see exercise below) and differentiable. Convexity of  $g$  now gives

$$g(t) \leq (1-t)g(0) + tg(1) \implies \frac{g(t) - g(0)}{t} \leq g(1) - g(0) \quad (44)$$

which means

$$g'(0) \leq g(1) - g(0). \quad (45)$$

Similarly we have

$$\frac{g(t) - g(1)}{t-1} \geq g(1) - g(0) \implies g'(1) \geq g(1) - g(0). \quad (46)$$

Thus

$$g'(1) \geq g'(0) \tag{47}$$

which translates back to:

$$Df(\mathbf{x})(\mathbf{x} - \mathbf{y}) \geq Df(\mathbf{y})(\mathbf{x} - \mathbf{y}) \tag{48}$$

and the proof ends. □

**Exercise 5.** Let  $f: \mathbb{R}^N \mapsto \mathbb{R}$ . Let  $\mathbf{x}_0, \mathbf{v} \in \mathbb{R}^N$  and define  $g: \mathbb{R} \mapsto \mathbb{R}$  by

$$g(t) = f(\mathbf{x}_0 + t \mathbf{v}). \tag{49}$$

- a) Prove that if  $f$  is convex then  $g$  is convex.
- b) If for every  $\mathbf{x}_0, \mathbf{v} \in \mathbb{R}^N$  the  $g$  as defined above is convex, can we conclude  $f$  is convex? Justify.